

# ScalAH22 (in conjunction w/SC22)

---



## Mixed-Precision Algorithm for Finding Selected Eigenvalues and Eigenvectors of Symmetric and Hermitian Matrices

Yaohung (Mike) Tsai<sup>1,4</sup>, Piotr Luszczek<sup>1</sup>, Jack Dongarra<sup>1,2,3</sup>

November 13, 2022

<sup>1</sup>University of Tennessee, <sup>2</sup>Oak Ridge National Laboratory, <sup>3</sup>University of Manchester

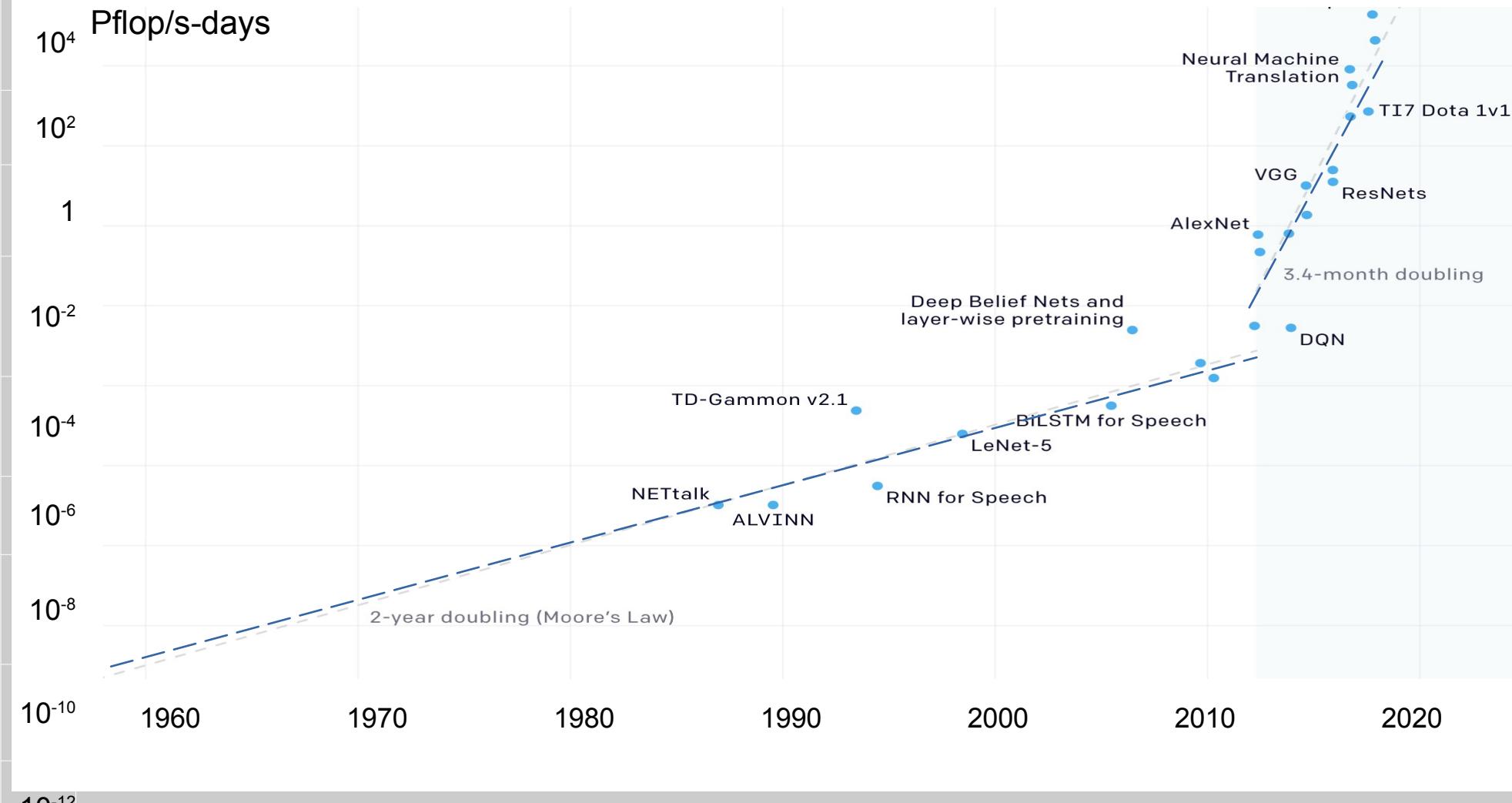
<sup>4</sup>Meta (formerly Facebook)

# The Landscape of Mixed-Precision Hardware

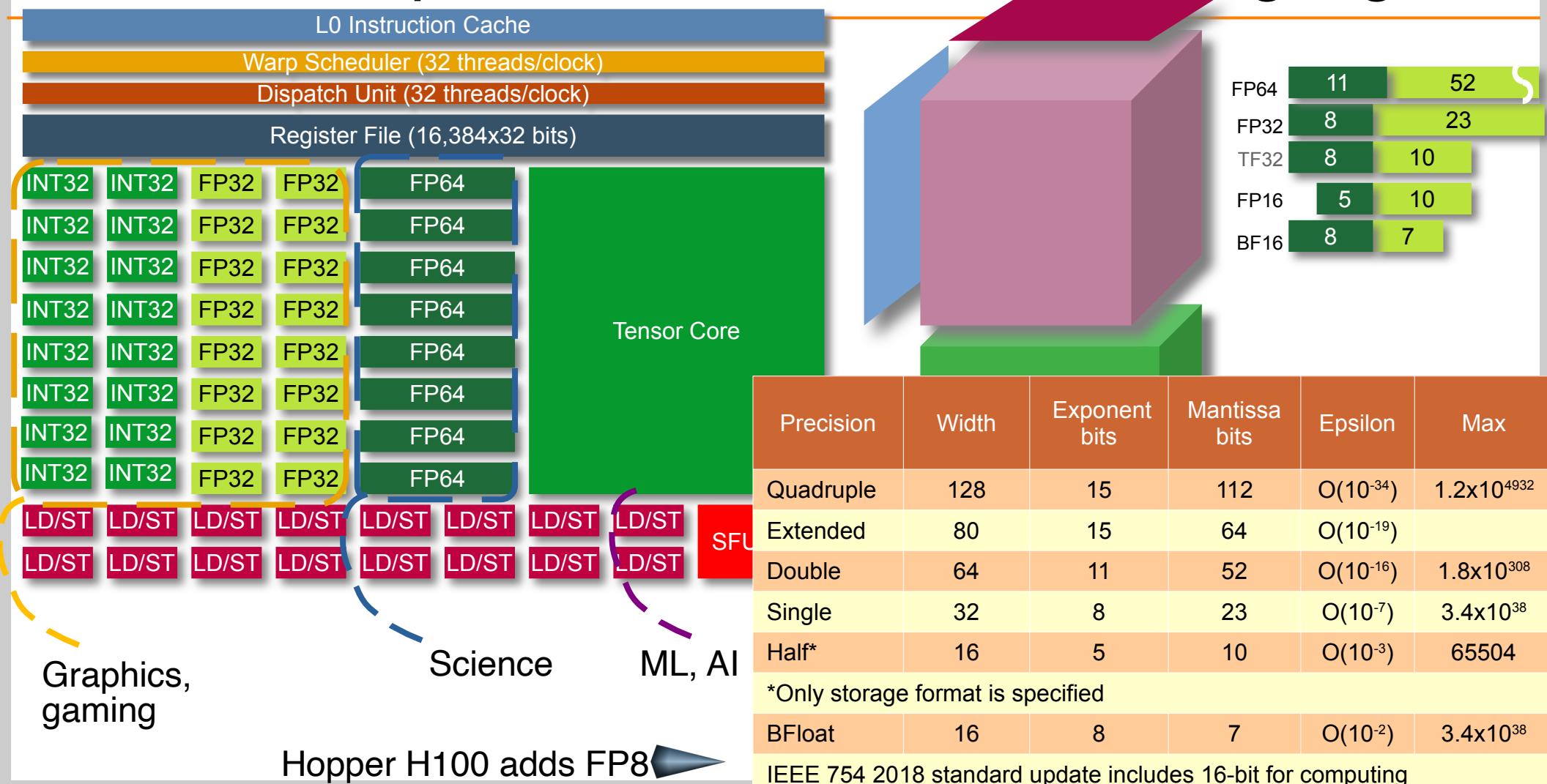
---

- Mixed-Precision Startups and their hardware
  - GraphCore
    - Colossus
  - Habana
    - Labs Gaudi
  - Cerebras
    - Wafer Scale Chip
  - Blaize
    - Graph Streaming Processor
  - Groq
    - Tensor Streaming Processor
  - SambaNova
    - Cardinal
  - Tenstorrent
    - Grayskull
- NVIDIA mixed-precision hardware
  - Pascal
    - FP16 units only
  - Volta
    - Tensor Cores and FP16
  - Turing
    - Tensor Cores and FP16
  - Ampere
    - Tensor Cores for FP16 and FP64
  - Hopper
    - Double-, single-, half-, and quarter-precision formats in scalar and tensor units

# Computational Needs According to OpenAI



# NVIDIA Ampere FP Hardware A100 Highlights



# Related Work for Eigenvalue Refinement

---

- Dongarra, J. J., Moler, C. B., & Wilkinson, J. H. (1983). *Improving the accuracy of computed eigenvalues and eigenvectors*. SIAM Journal on Numerical Analysis, 20(1), 23-45.
- Prikopa, K. E., & Gansterer, W. N. (2013). *On mixed precision iterative refinement for eigenvalue problems*. Procedia Computer Science, 18, 2647-2650.
- Ogita, T., & Aishima, K. (2018). *Iterative refinement for symmetric eigenvalue decomposition*. Japan Journal of Industrial and Applied Mathematics, 35(3), 1007-1035.

## Algorithm 1 SICE algorithm (Dongarra, Moler, Wilkinson)

1: **Input:** Matrix  $A \in \mathbb{R}^{n \times n}$ , An approximate eigenvalue  $\lambda$  and the corresponding eigenvector  $x$ .  $I_{\max}$  denotes the maximum number of iterations.

2: **Output:** Refined eigenvalue  $\lambda$  and its eigenvector  $x$ .

3: **function**  $[\lambda, x] \leftarrow \text{SICE}(A, \lambda, x, I_{\max})$

4:    $[Q, U] \leftarrow \text{schur}(A)$  ► obtain Schur decomposition  
 $A = QUQ^T, QQ^T = I.$

5:    $[m, s] \leftarrow \max(\text{abs}(x)); x \leftarrow x/m$  ► Normalizing  $x$  so that  $\|x\|_\infty = s_x = 1$ .

6:   **for**  $i$  in  $1 : I_{\max}$  **do**

7:      $r \leftarrow \lambda x - Ax$

8:      $c \leftarrow -x - a_{\lambda s}$

9:      $d \leftarrow Q^T c$

10:     $f^T \leftarrow Q(s, :) = e_s^T Q$  ►  $s$ -th row of  $Q$ .

11:     $\bar{U}_\lambda \leftarrow Q_1(U - \lambda I); \bar{d} \leftarrow Q_1 d = \|d\|_2 e_1$  ► Givens rotations  $Q_1$  from  $Q_1 d = (P_2 P_3 \dots P_n) d = \gamma e_1$  where  $\gamma = \|d\|_2$

12:     $\bar{U}_\lambda \leftarrow \bar{U}_\lambda + \bar{d}(1)f^T$

13:     $\bar{U}_\lambda \leftarrow Q_2 \bar{U}_\lambda$  ► Givens rotations  $Q_2$  to introduce upper triangular form.

14:    Solve the triangular system  $\bar{U}_\lambda z = Q_2 Q_1 Q^T r$

15:     $y \leftarrow Qy$

16:     $\lambda \leftarrow \lambda + y(s)$  ► Update eigenvalue.

17:     $y(s) \leftarrow 0$  ► Set  $y(s)$  to 0.

18:     $x \leftarrow x + y$  ► Update eigenvector.

19:    **if** desired accuracy is reached **then**

20:     **break**

21:    **end if**

22: **end for**

23: **end function**

```
[Q, T]=schur(single(A)); Q = V; T = D;
result=zeros(1,niter+1); normA=norm(A,inf);
for i=1:niter
    r=lambda*x-A*x; %residual vector
    error(i)=norm(r,inf)/normA/norm(x,inf);
    disp(['Norm of residual vector before iteration ', num2str(i), ' : ',
    num2str(norm(r))]);
    a_lambda_s = A(:,s); % eq (3.14)
    a_lambda_s(s) = a_lambda_s(s) - lambda;
    c=-x-a_lambda_s;
    d=Q'*c; % eq (3.15)
    fT=Q(s,:);
    %
    % Explicitly form matrix B: B=A-lambda*eye(n); B(:,s)=B(:,s)+c; y=B\z;
    %
    T_lambda=T-lambda*eye(n);
    for k=n:-1:2 % Apply Q_1
        [P, z] = planerot(d(k-1:k)); % givens rotation
        %T_lambda(k-1:k,:) = P*T_lambda(k-1:k,:);
        T_lambda(k-1:k,k-1:n) = P*T_lambda(k-1:k,k-1:n);
        d(k-1:k) = z;
        rhs(k-1:k) = P * rhs(k-1:k);
    end
    T_lambda(1,:) = T_lambda(1,:) + d(1)*fT;
    % Apply Q_2
    for k=1:n-1
        [P, z] = planerot(T_lambda(k:k+1,k));
        %T_lambda(k:k+1,:) = P*T_lambda(k:k+1,:);
        T_lambda(k:k+1,k:n) = P*T_lambda(k:k+1,k:n);
        rhs(k:k+1) = P*rhs(k:k+1);
        T_lambda(k:k+1,k)=z;
        %T_lambda(k+1,k)=0; % forcing real zero
    end
    y=T_lambda\rhs; % T_lambda is upper triangular
    %y=(T-(lambda)*eye(n)+d*fT)\rhs;
    %y=gmres(T_lambda,rhs,5,1e-12,15,T_lambda);
    y=Q*y;
    new_x=x+y; % new eigenvector
    new_x(s)=x(s); % restore x(s) since y_s = 0
    disp(['Computed eigenvalue offset at iteration ', num2str(i), ' : ',
    num2str(y(s))]);
    lambda=lambda+y(s); % new eigenvalue
    % convergence check: break if 2*(previous correction) < current correction
    %if(i>=2 && abs(y(s))>2*abs(last_u))
    last_u = y(s);
    x=new_x;
    end
```

# Mixed-Precision Eigensolver Goals

---

$$Ax = \lambda x \text{ where } A = A^T \vee A = A^H$$

- Looking into iterative refinement for real symmetric and complex Hermitian eigenvalue problems.
- Apply mixed-precision to improve the performance.
- Incorporate with latest two-stages eigensolvers.
- Targeting the case that only part of the eigenpairs are required.

---

**Algorithm 2** SICE-SM algorithm: SICE algorithm with Sherman–Morrison formula

---

- 1: **Input:** Matrix  $A = A^\top \in \mathbb{R}^{n \times n}$ . An approximate eigenvalue  $\lambda$  and the corresponding eigenvector  $x$ .  $I_{\max}$  denotes the maximum number of iterations.
- 2: **Output:** Refined eigenvalue  $\lambda$  and eigenvector  $x$ .
- 3: **function**  $[\lambda, x] \leftarrow \text{SICE\_SM}(A, \lambda, x, I_{\max})$
- 4:    $[Q, T] \leftarrow \text{tridiag}(A)$                $\triangleright$  Tridiagonalization  $A = QTQ^\top$ ,  
     $QQ^\top = I$ .
- 5:    $[m, s] \leftarrow \max(\text{abs}(x)); x \leftarrow x/m$        $\triangleright$  Normalization of  $x$  so  
    that  $\|x\|_\infty = s_x = 1$ .
- 6:   **for**  $i$  in  $1 : I_{\max}$  **do**
- 7:      $r \leftarrow \lambda x - Ax$
- 8:      $c \leftarrow -x - a_{\lambda s}$
- 9:      $d \leftarrow Q^\top c$
- 10:     $f^\top \leftarrow Q(s, :) = e_s^\top Q$                        $\triangleright$   $s$ -th row of  $Q$ .
- 11:     $rhs \leftarrow Q^\top r$
- 12:     $u \leftarrow (T - \lambda I)^{-1}d$
- 13:     $v \leftarrow (T - \lambda I)^{-1}rhs$
- 14:     $y \leftarrow v - \frac{f^\top v}{1 + f^\top u} u$        $\triangleright$  Sherman–Morrison formula
- 15:     $y \leftarrow Qy$
- 16:     $\lambda \leftarrow \lambda + y(s)$                        $\triangleright$  Update eigenvalue.
- 17:    **if**  $i \neq 1$  **then**
- 18:       $y(s) \leftarrow 0$                        $\triangleright$  Set  $y(s)$  to 0.
- 19:       $x \leftarrow x + y$                        $\triangleright$  Update eigenvector.
- 20:    **end if**
- 21:    **if** desired accuracy reached **then**
- 22:      **break**
- 23:    **end if**
- 24:   **end for**
- 25: **end function**

# Background — SICE algorithm

---

- General matrix  $A$  with approximated eigenvalue  $\lambda$  and eigenvector  $x$ .
- Solving a rank one updated system to obtain the correction to both eigenvalue  $y_s$  and eigenvector  $y$ :

$$(A - \lambda I + c e_s^T) y = \lambda x - Ax$$

- Solving with Schur decomposition  $A = QUQ^H$  and two series of Givens rotations.

# New SICE-SM Algorithm

---

- Looking at symmetric matrices first:  $A = A^T$
- The rank one updated system under tridiagonalization  $A = QTQ^T$

$$Q(T - \lambda I + Q^T c e_s^T Q) Q^T y = \lambda x - Ax$$

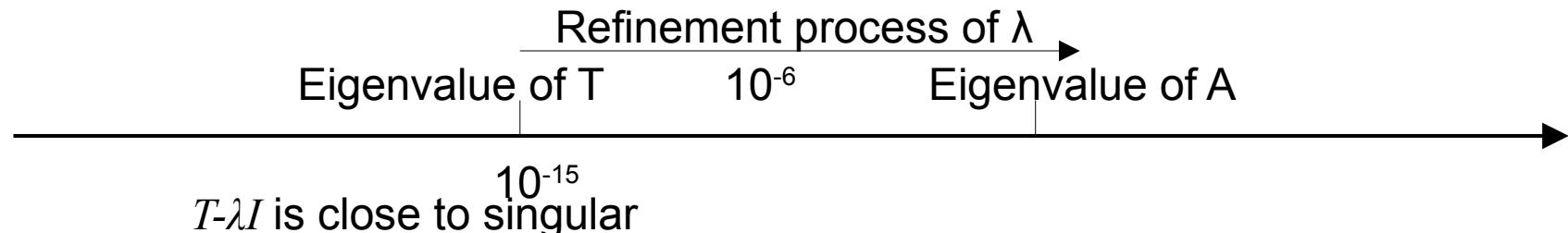
- Using Sherman-Morrison formula

$$(A - uv^T)^{-1} = A^{-1} - \frac{A^{-1} uv^T A^{-1}}{1 + v^T A^{-1} u}$$

- How close is  $T - \lambda I$  to being singular? In which precision?

# Mixed Precision Eigensolver Process

- Consider single and double precision.
  - Tridiagonalization  $A = QTQ^T$  in single precision.
- Find the eigenvalues  $\Lambda$  and eigenvectors  $V$  of the tridiagonal system  $T$  in double precision.
- Back transform to obtain the eigenvectors of  $A$ :  $V \leftarrow QV$
- Iterative refinement with SICE-SM



### Algorithm 3 Blocked SICE-SM algorithm

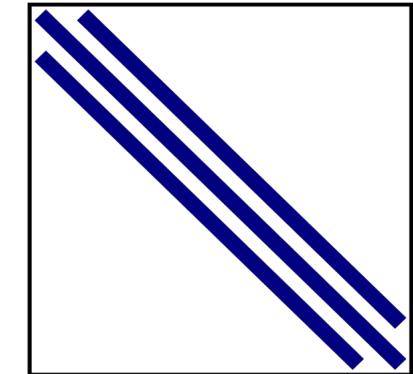
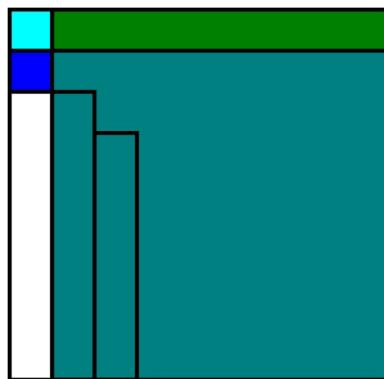
1: **Input:**  $A = A^T \in \mathbb{R}^{n \times n}$ , initial eigenvectors  
 $X = [x_1 | x_2 | \dots | x_\ell] \in \mathbb{R}^{n \times \ell}$  and the corresponding initial  
eigenvalues  $\Lambda = (\lambda_1, \lambda_2, \dots, \lambda_\ell)^T \in \mathbb{R}^\ell$ .  $I_{\max}$  denotes the  
maximum number of iterations.  
2: **Output:** Refined eigenvectors  $X$  and refined eigenvalues  $\Lambda$ .  
3: **function**  $[X, \Lambda] \leftarrow \text{SICE\_SM\_BLK}(A, X, \Lambda, I_{\max})$   
4:    $[Q, T] \leftarrow \text{tridiag}(A)$                $\triangleright$  Tridiagonalization  $A = QTQ^T$ ,  
     $QQ^T = I$ .  
5:   **for**  $i$  in  $1 : I_{\max}$  **do**  
6:      $s \leftarrow i$   
7:      $R \leftarrow X \times \text{diag\_matrix}(\Lambda) - A \times X$        $\triangleright$  Residual vectors  
    need higher precision.  
8:     **for**  $j$  in  $1 : \ell$  **do**  
9:        $c_j \leftarrow -x_j - A(:, s)$   
10:      **end for**  
11:      Compose matrix  $C = [c_1 | c_2 | \dots | c_\ell]$  from column vectors  
     $c_j$   
12:       $C(s, :) \leftarrow C(s, :) + \Lambda^T$   
13:       $D = [d_1 | d_2 | \dots | d_\ell] \leftarrow Q^T \times C$   $\triangleright$  Can be in lower precision.  
14:       $RHS = [rhs_1 | rhs_2 | \dots | rhs_\ell] \leftarrow Q^T \times R$   $\triangleright$  Can be in lower  
    precision.  
15:       $f \leftarrow Q(s, :)$                                $\triangleright$   $s$ -th row of  $Q$ .  
16:      **for**  $j$  in  $1 : \ell$  **do**  
17:        $u_i \leftarrow (T - \lambda J)^{-1} d_i$   
18:        $v_i \leftarrow (T - \lambda J)^{-1} rhs_i$   
19:        $y_i \leftarrow v_i - \frac{f^T u_i}{1 + f^T u_i} u_i$                $\triangleright$  Sherman-Morrison  
20:      **end for**

21:       $y_j$                                       Compose matrix  $Y = [y_1 | y_2 | \dots | y_\ell]$  from correction vectors  
22:       $Y \leftarrow Q \times Y$   
23:       $\Lambda \leftarrow \Lambda + Y(s, :)^T$                        $\triangleright$  Update eigenvalues.  
24:      **if**  $i \neq 1$  **then**  
25:        $Y(s, :) \leftarrow 0$                                $\triangleright$  Set  $y_i(s)$  to 0.  
26:        $X \leftarrow X + Y$                                $\triangleright$  Update eigenvectors.  
27:       Normalize eigenvectors  $x_i$  in  $X$ .  
28:      **end if**  
29:      **if** desired accuracy reached **then**  
30:       **break**  
31:      **end if**  
32:      **end for**  
33:       $X \leftarrow X + \frac{1}{2} X(I - X^T X)$                        $\triangleright$  Orthogonalization.  
34: **end function**

# Background: Two-Stage Tridiagonalization

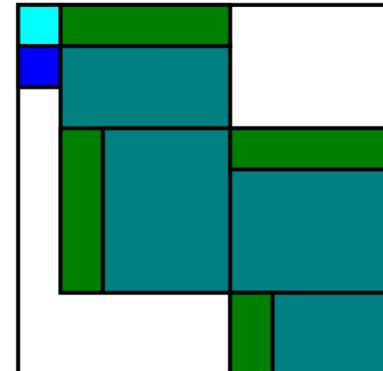
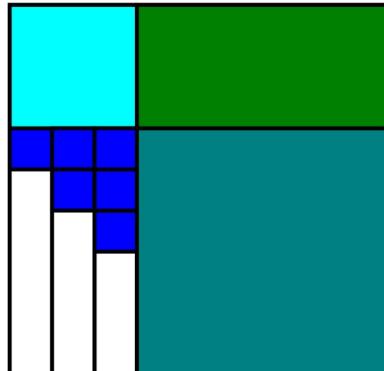
One stage (classic)

- Householder transformations
- Bound by memory speed



Two stage (post SBR)

- First stage: symmetric to band
- Second stage: band to tridiagonal



QR an LQ

Bulge chasing

# Blocked SICE-SM Algorithm

- Refine multiple eigenvalues and eigenvectors simultaneously
- Performance of  $n \times n$  matrix times  $n \times m$  aggregated vectors on NVIDIA V100

Matrix size	Number of eigenvectors	Time (ms)	Performance (Tflop/s)
20000	1	3.76	0.212
20000	8	3.79	1.688
20000	32	6.48	3.949
20000	128	13.57	7.544

- Reorthogonalize with one Newton–Schulz iteration at the end:  
$$X' = X + \frac{1}{2} X(I - X^T X)^{-1}$$

# Implementation Details

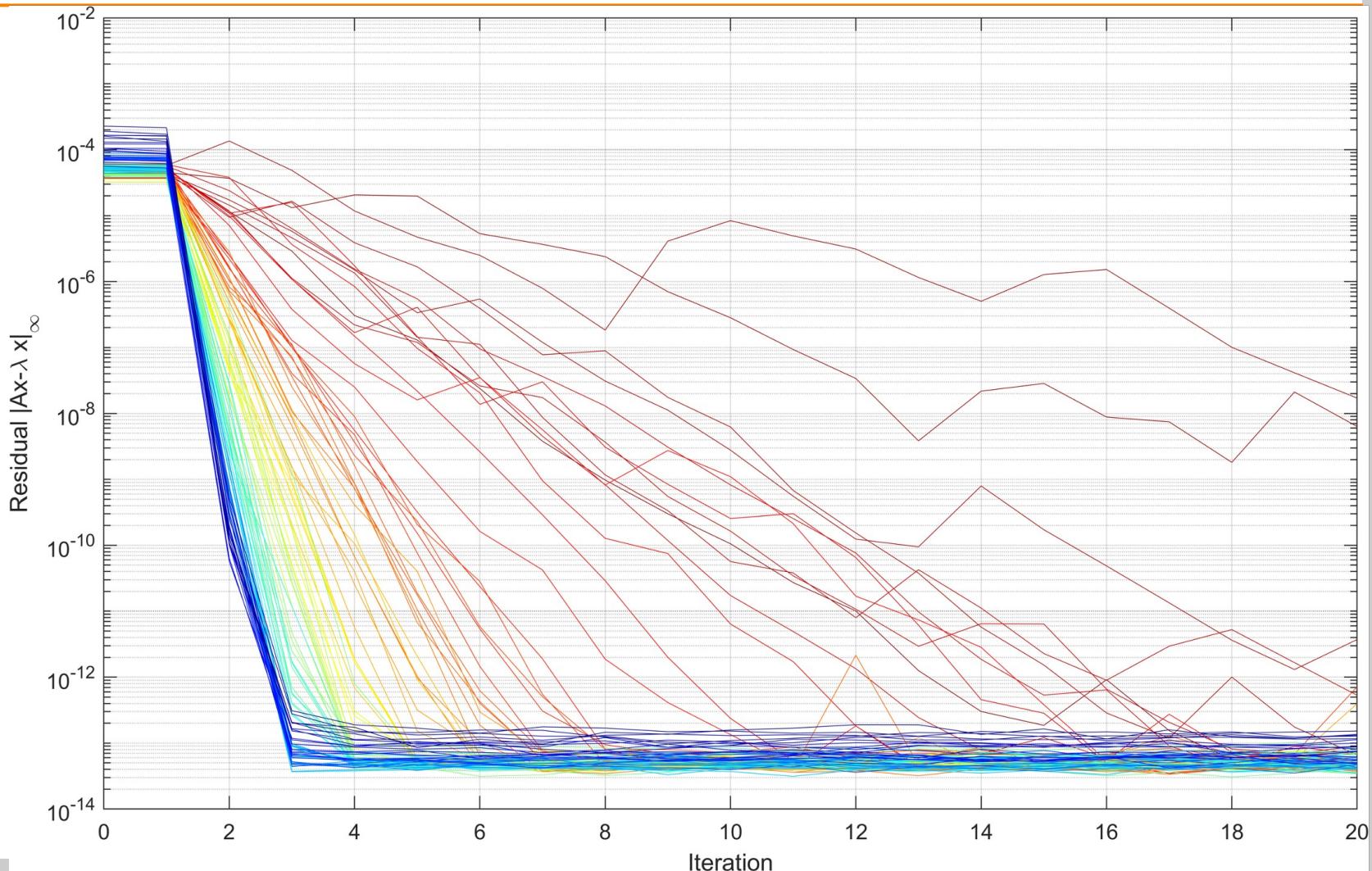
- MAGMA numerical library
  - Targeting heterogeneous systems (CPU+GPU)
- Reordering the back transformations and overlapping with CPU work  
*(not to scale)*



- Batched tridiagonal solver on GPU

# Numerical Convergence Results

- $\text{size}(A)=100$
- $\text{cond}(A)=10^7$
- geometrically distributed eigenvalues from 1 to  $10^{-7}$



# Performance Results on NVIDIA Volta V100

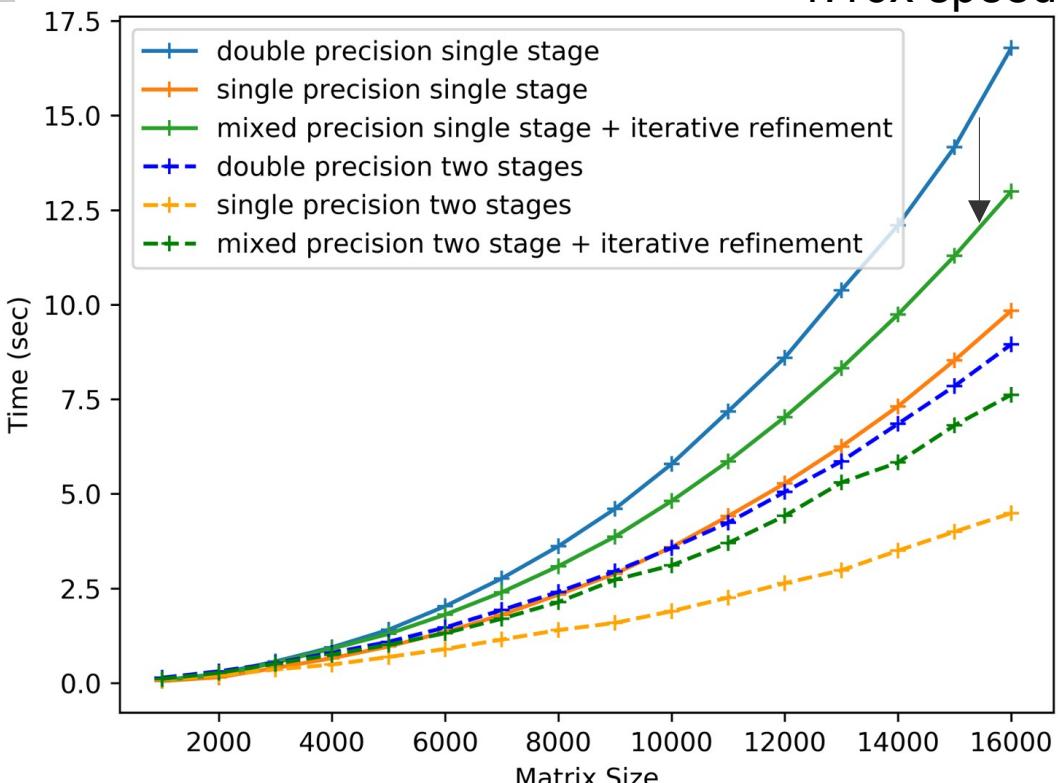
CPU: Intel Xeon E5-2650 v3 (Haswell 10C20T), GPU: NVIDIA Tesla Volta V100.

Intel compilers and MKL, CUDA 11, MAGMA 2.5.4.

Requesting largest 32 eigenpairs.

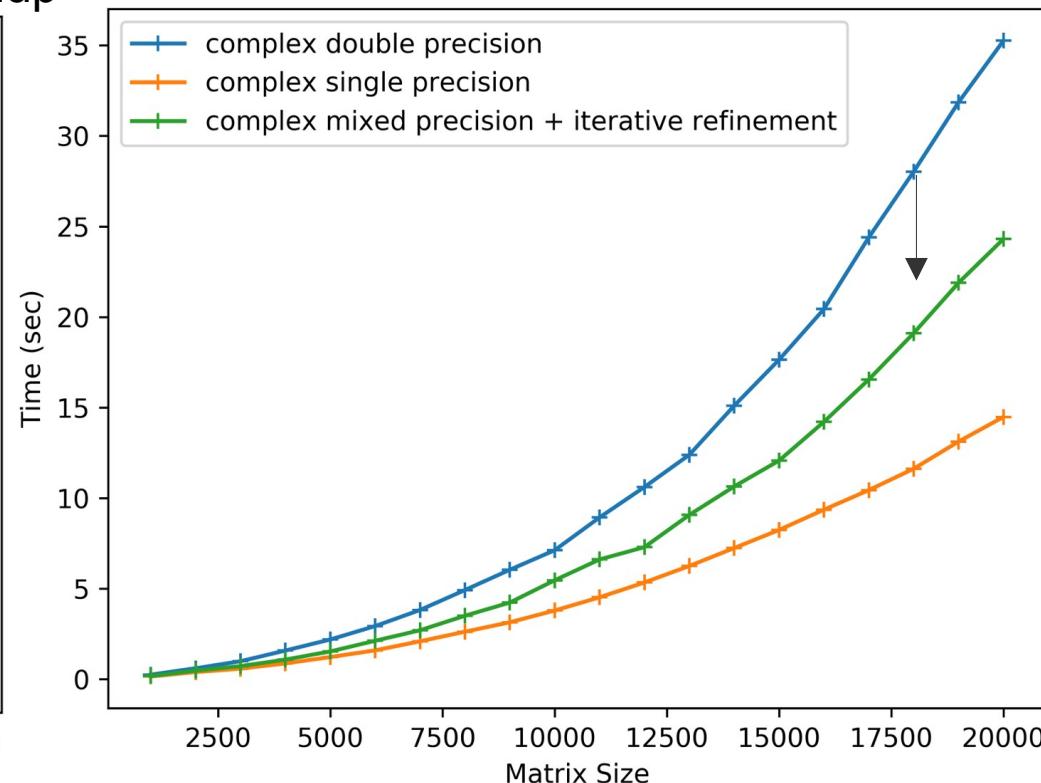
IR

1.16x speedup



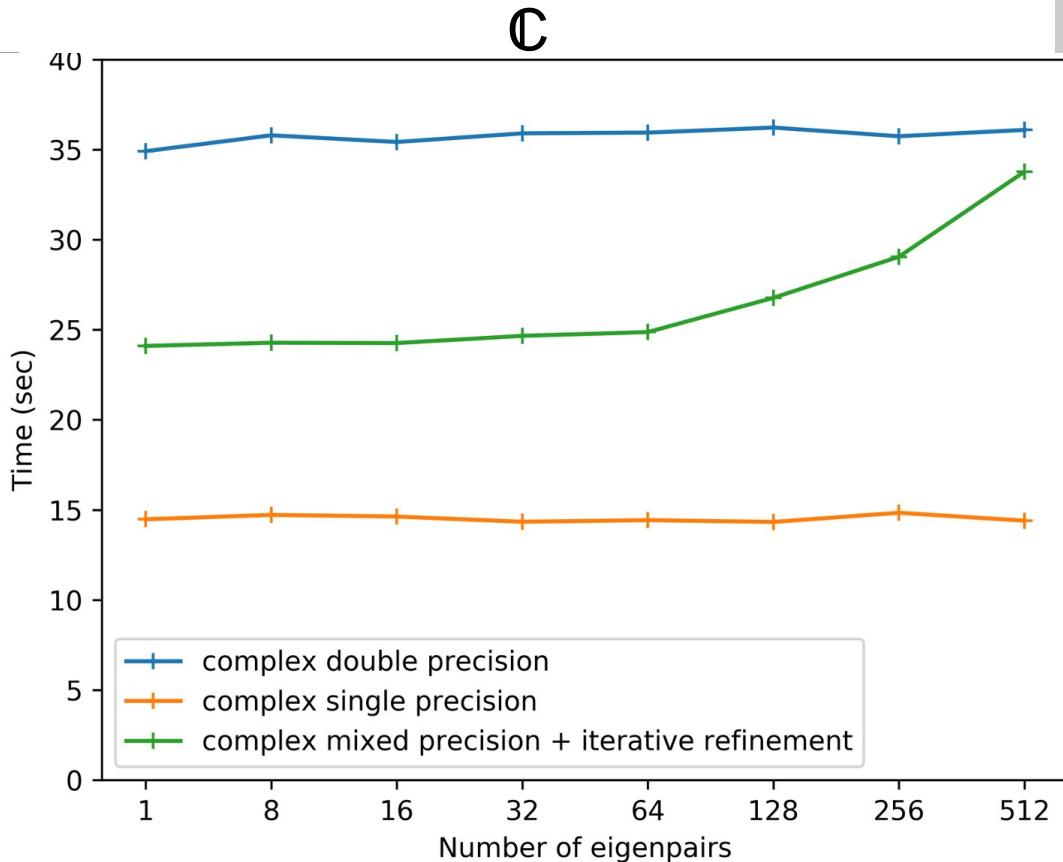
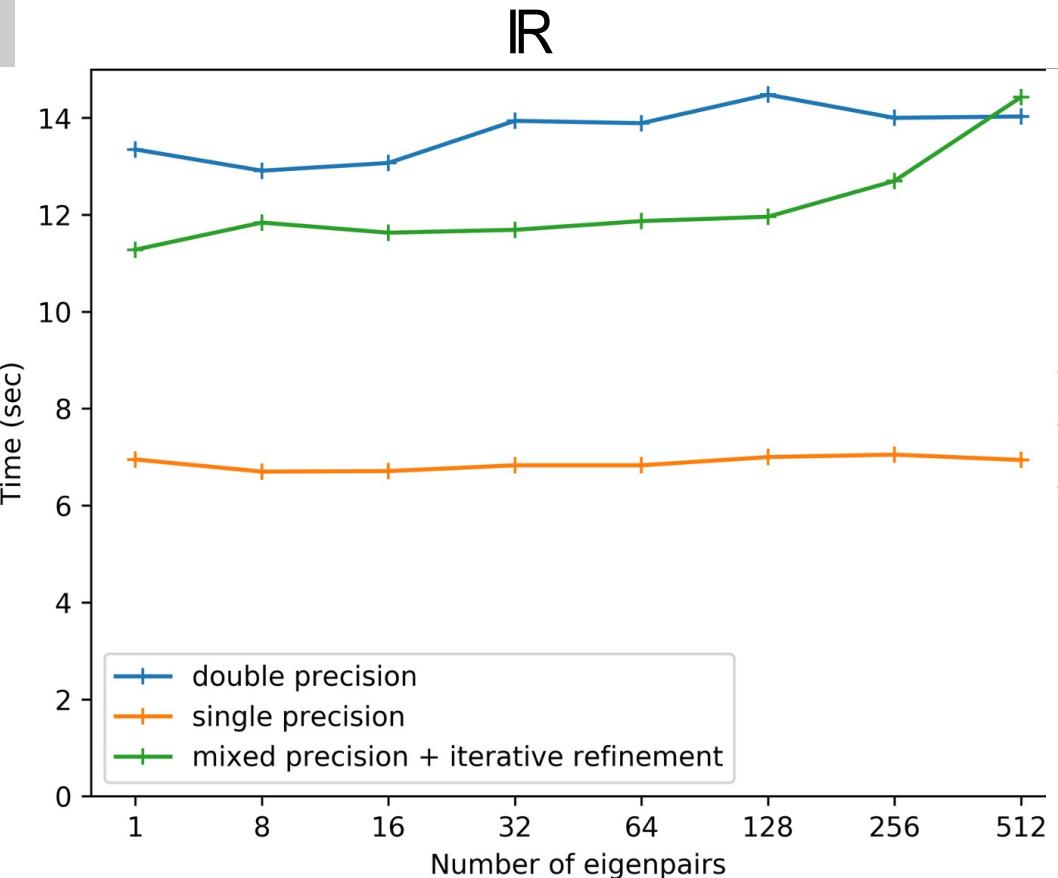
C

1.45x speedup



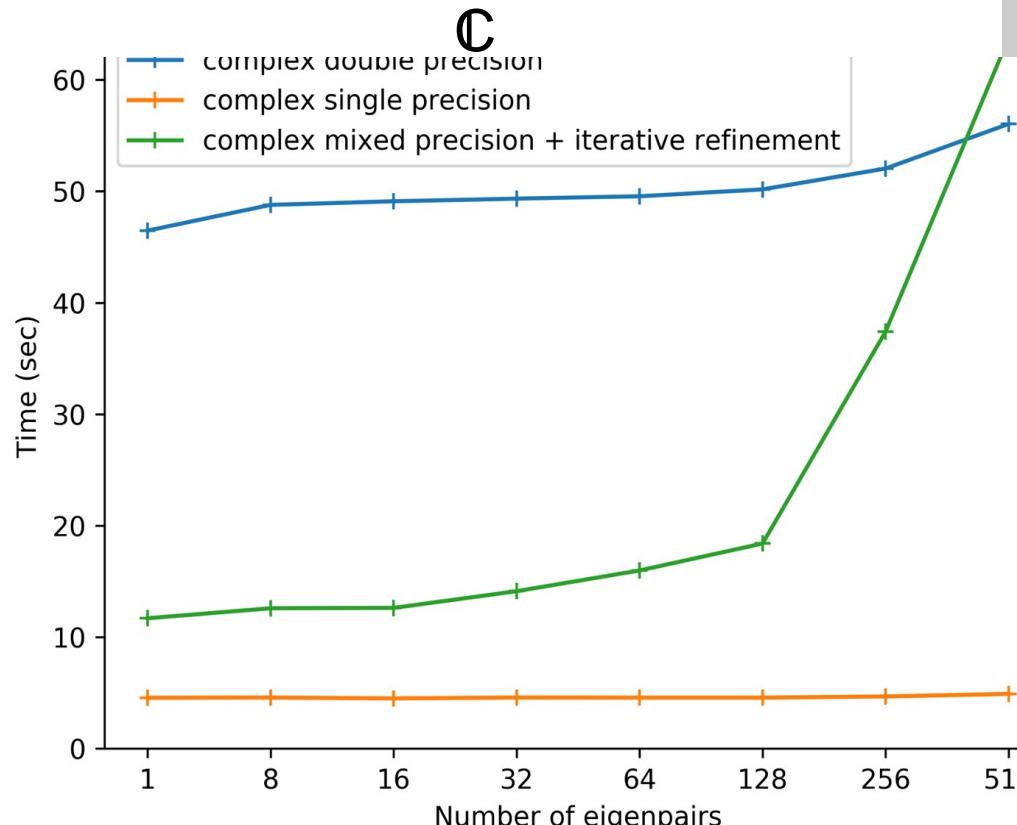
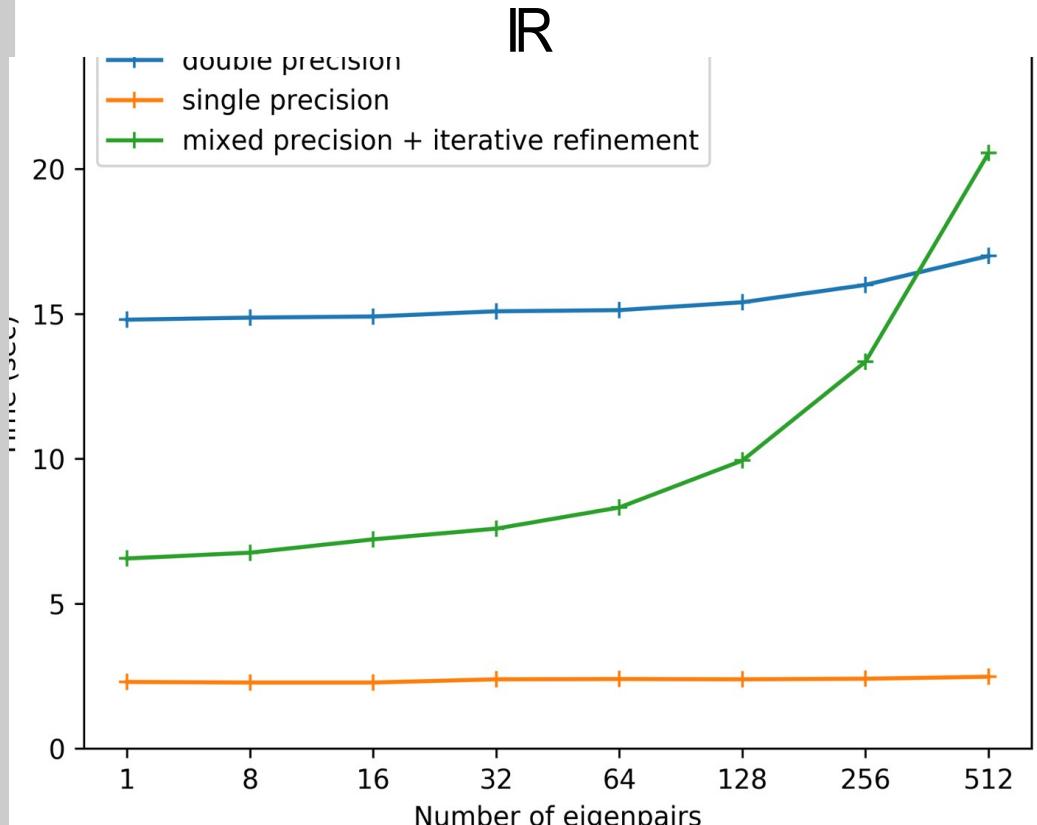
# Time for Refining Different Number of Eigenpairs

Matrix size N=20000, GPU: NVIDIA Tesla Volta V100



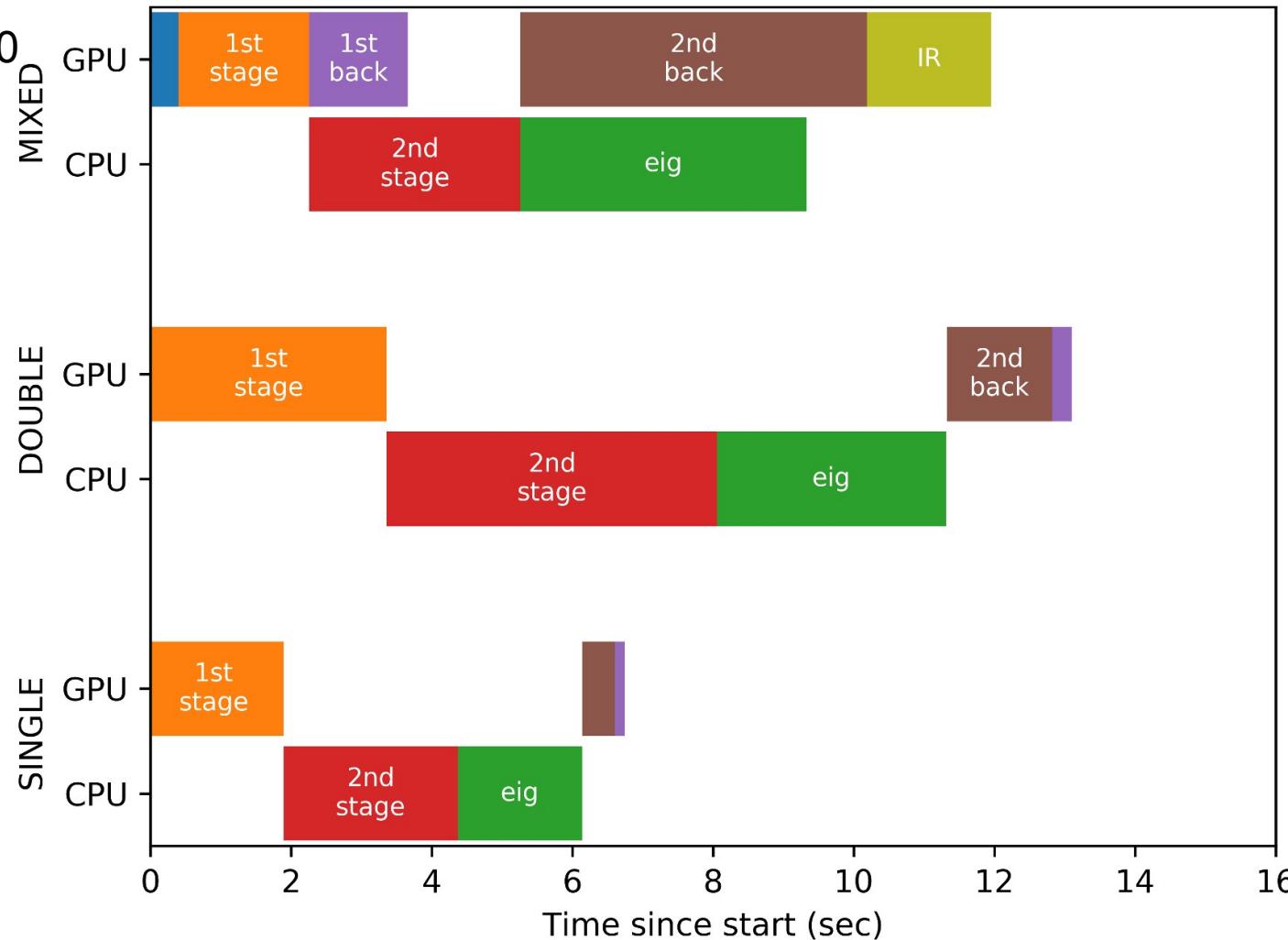
# Time for Refining Different Number of Eigenpairs

Matrix size N=20000, GPU: NVIDIA Tesla Pascal GTX1060



# Runtime Trace

Real symmetric N=20000  
32 eigenpairs



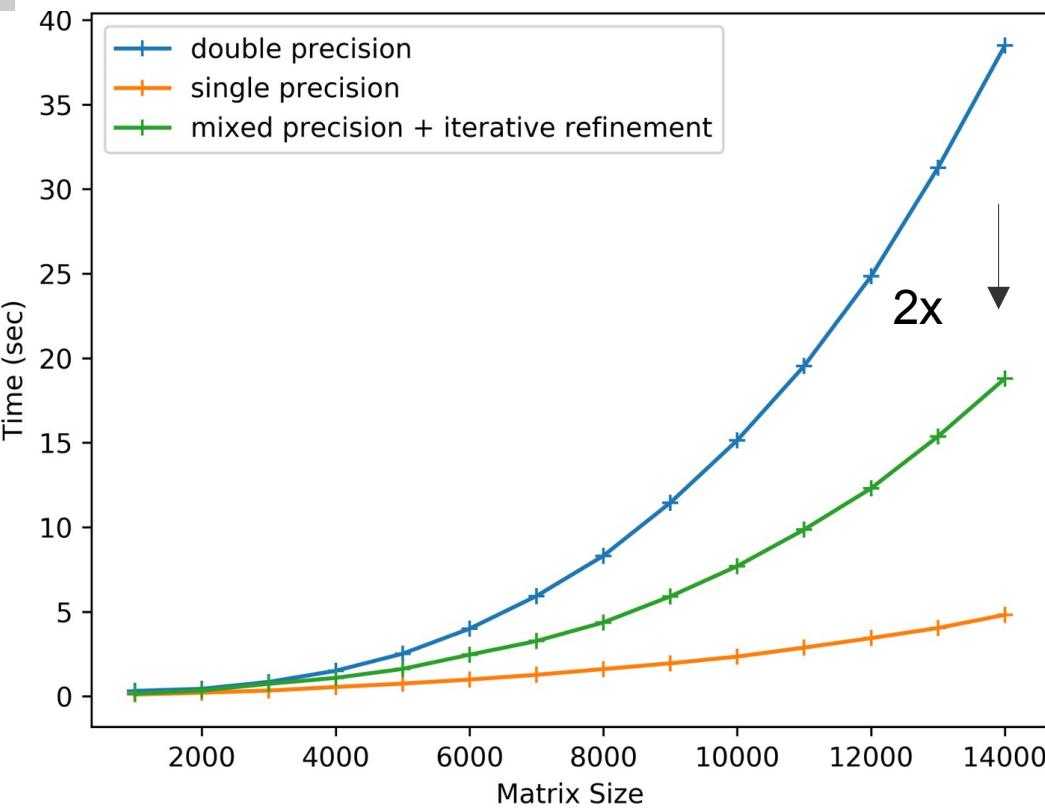
# Performance Results on NVIDIA GTX1060

GPU: NVIDIA Pascal GTX1060

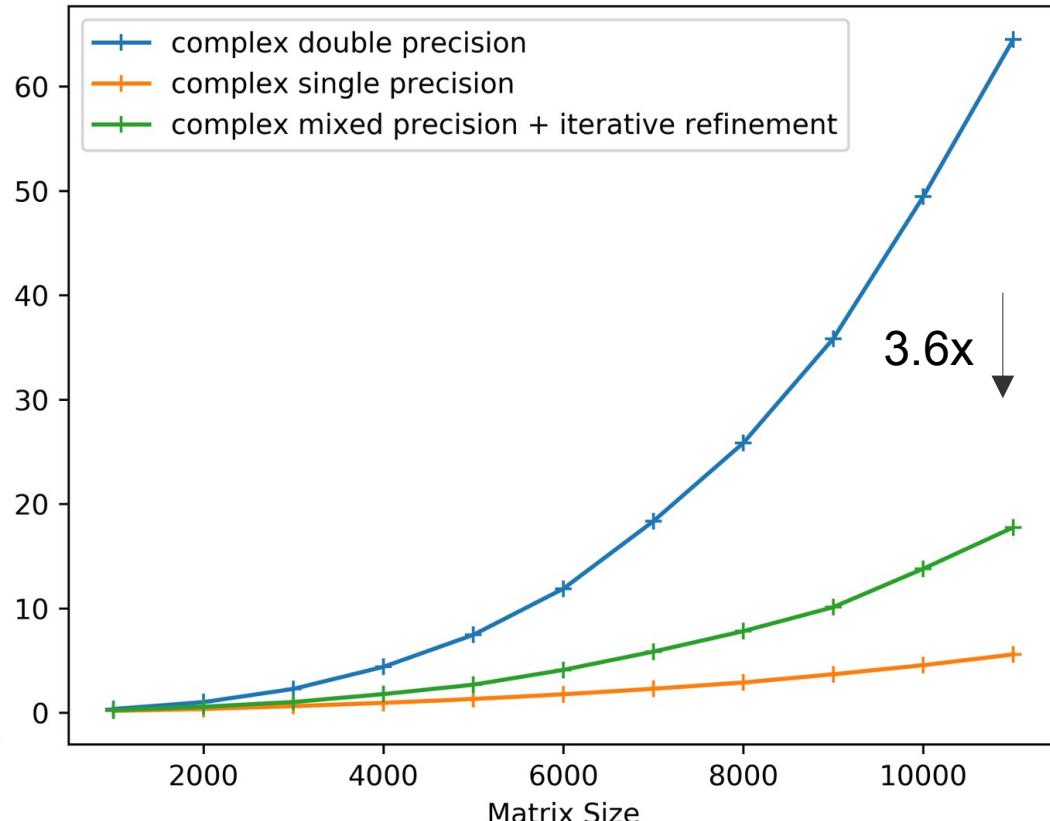
Peak single 4.375 TFLOP/s, double 136.7 GFLOP/s (1:32).

Requesting largest 32 eigenpairs.

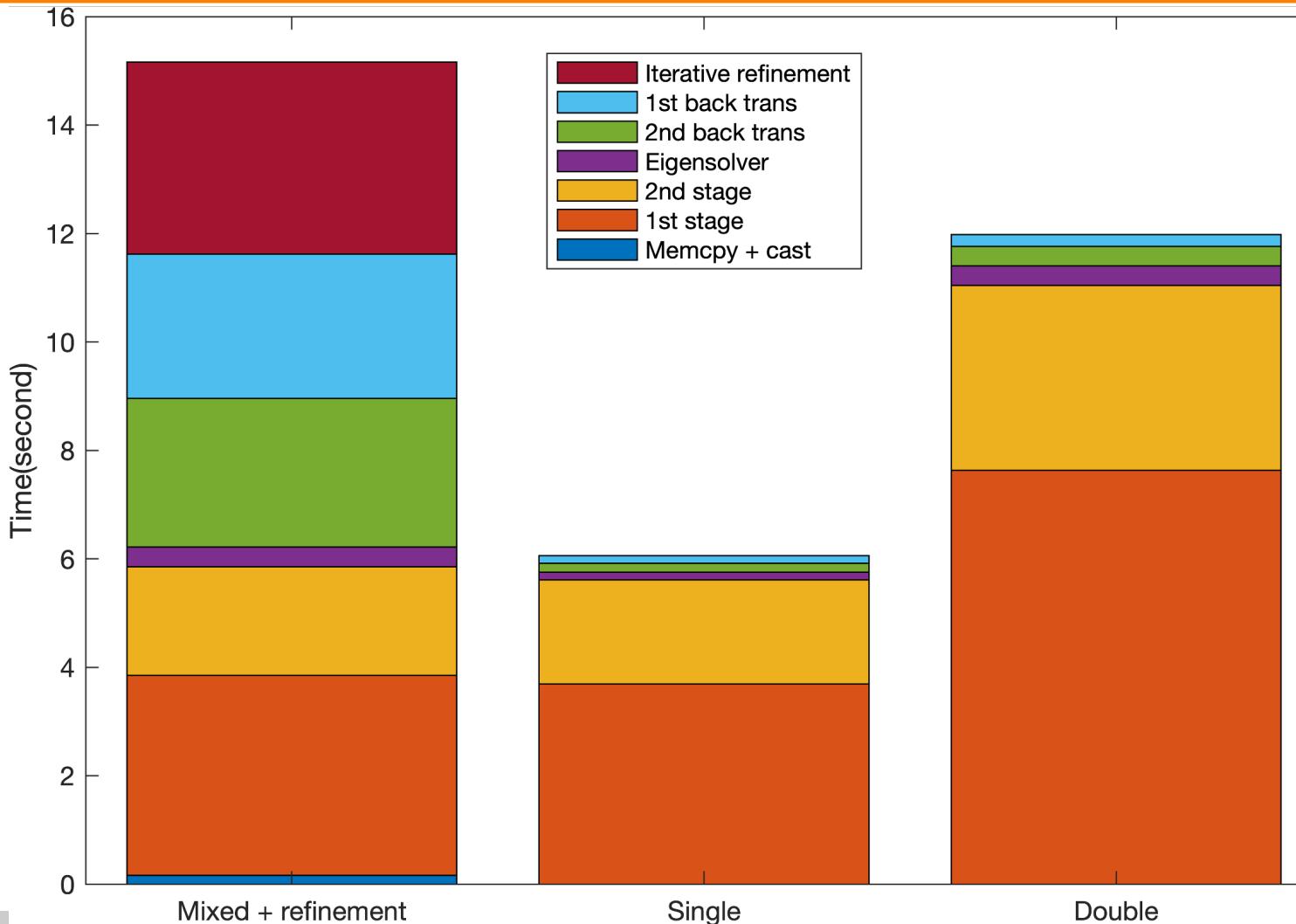
$\mathbb{R}$



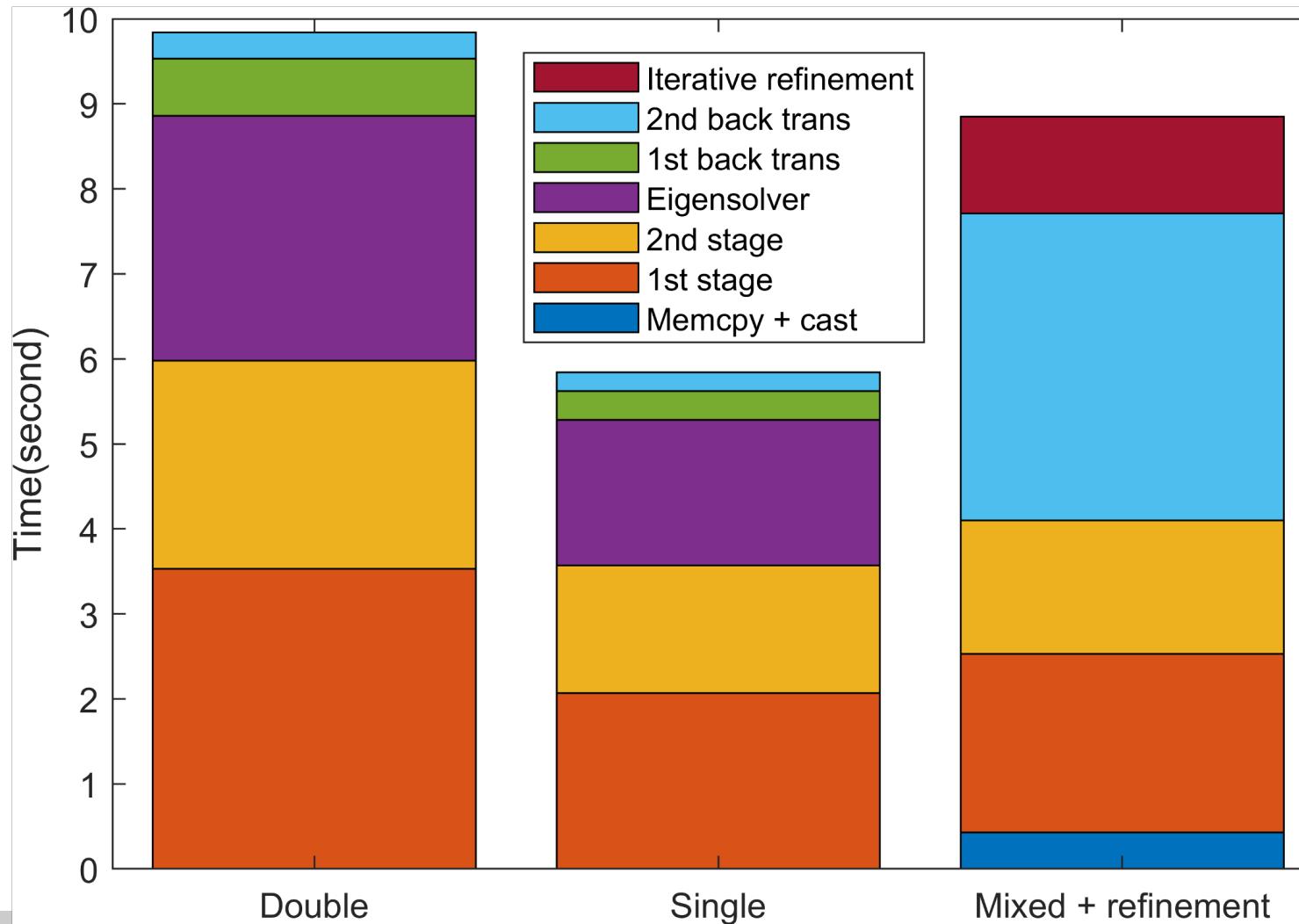
$\mathbb{C}$



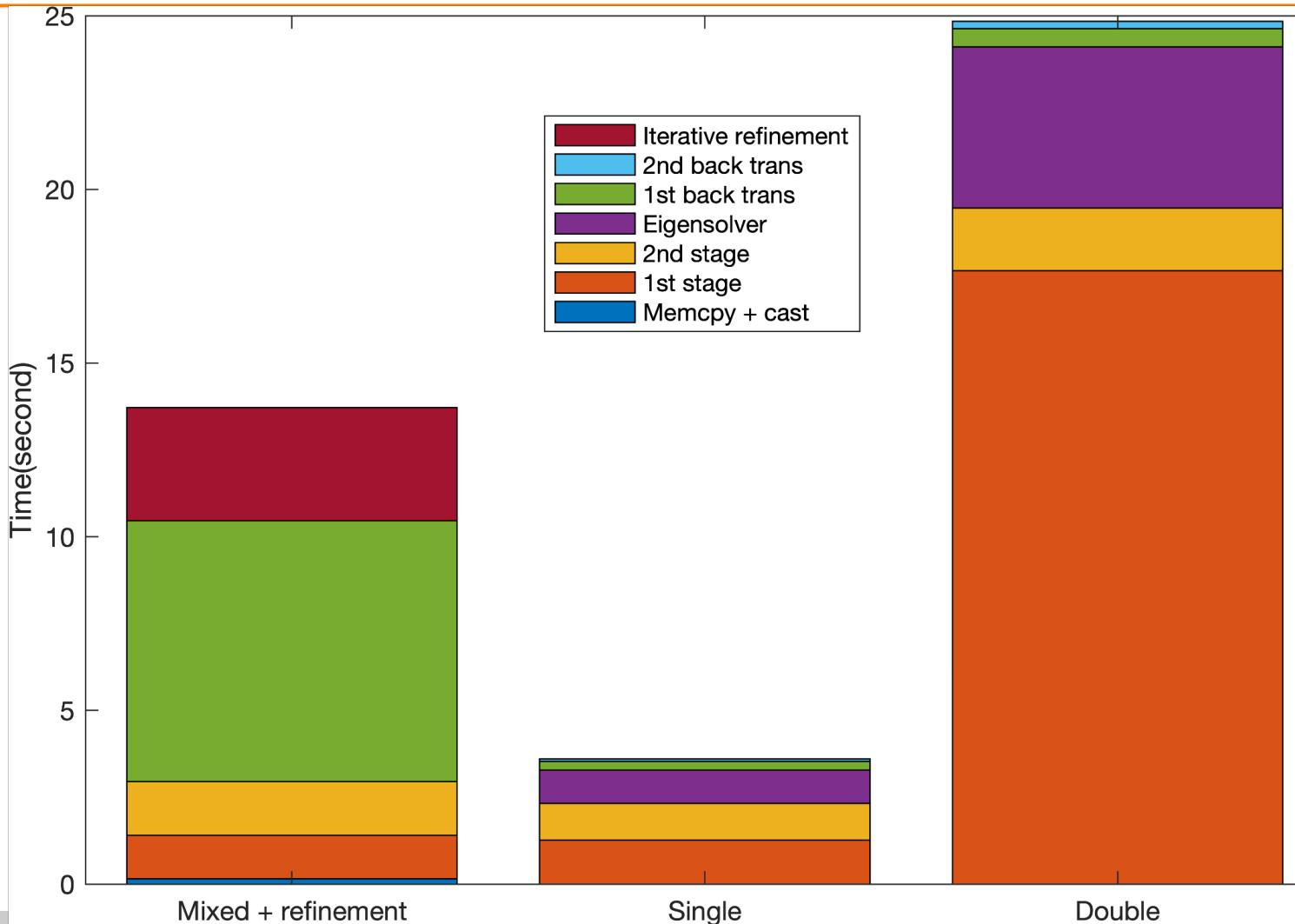
# Time Breakdown on Multicore CPU



# Time Breakdown on NVIDIA Volta V100



# Time Breakdown on NVIDIA Pascal GTX1060



# Mixed-Precision Eigensolver Summary

---

- Methods' summary
  - Blocked version of SICE algorithm that uses Sherman-Morrison formula
  - There are performance improvements over the state-of-the-art two-stage algorithms
- Possible future extensions
  - Convergence analysis
    - Dealing with near-singularity of  $T - \lambda I$
    - $Q^T A Q$  isn't tri-diagonal while  $Q^T_{32} A Q_{32}$  is
      - Also: orthogonality of  $Q$  depends on precision
  - Dealing with clustered eigenvalues
  - Implement more efficiently the second stage that forms  $Q$
  - Use more of the custom low-precision formats
  - Implementation for distributed memory systems

# ICL @ SC22

UNIVERSITY OF TENNESSEE **BOOTH #613**



@ICL\_UTK ON TWITTER

FOLLOW US AT <https://icl.utk.edu/sc22>